

A stochastic approach for quantifying immigrant integration: the Spanish test case.

Elena Agliari,^{1,2} Adriano Barra,^{1,2} Pierluigi Contucci,³ Rickard Sandell,⁴ and Cecilia Vernia⁵

¹*Dipartimento di Fisica, Sapienza Università di Roma, Piazzale Aldo Moro 2, 00185, Roma, Italy*

²*INdAM, Gruppo Collegato dell'Università di Roma "Tor Vergata",
Dipartimento di Matematica, via della Ricerca Scientifica 1, 00133, Roma, Italy*

³*Dipartimento di Matematica, Università degli Studi di Bologna,
piazza Porta San Donato 2, 00124, Bologna, Italy*

⁴*Departamento de Ciencias Sociales, Universidad de Carlos III de Madrid,
Avenida de la Universidad 30, 28911, Madrid, Spain*

⁵*Dipartimento di Matematica, Università degli Studi di Modena e Reggio, Corso Canal Grande 64 01234, Modena, Italy*

(Dated: May 6, 2014)

We apply stochastic process theory to the analysis of immigrant integration. Using a unique and detailed data set from Spain, we study the relationship between local immigrant density and two social and two economic immigration quantifiers for the period 1999–2010. As opposed to the classic time-series approach, by letting immigrant density play the role of "time", and the quantifier the role of "space" it become possible to analyze the behavior of the quantifiers by means of continuous time random walks. Two classes of results are obtained. First we show that social integration quantifiers evolve following pure diffusion law, while the evolution of economic quantifiers exhibit ballistic dynamics. Second we make predictions of best and worst case scenarios taking into account large local fluctuations. Our stochastic process approach to integration lends itself to interesting forecasting scenarios which, in the hands of policy makers, have the potential to improve political responses to integration problems. For instance, estimating the standard first-passage time and maximum-span walk reveals local differences in integration performance for different immigration scenarios. Thus, by recognizing the importance of local fluctuations around national means, this research constitutes an important tool to assess the impact of immigration phenomena on municipal budgets and to set up solid multi-ethnic plans at the municipal level as immigration pressure build.

PACS numbers: 89.65.Ef, 89.65.-s, 05.40.-a, 05.70.Fh

I. INTRODUCTION

A particular political challenge of growing immigration is immigrant integration. It is considered a necessity for minimizing frictions and confrontation between immigrants and natives in the host community, as well as a precondition for a competitive and sustainable economy[1]. In response to the recent rapid growth in the number of immigrants throughout many major regions in the world, the need for political intervention targeting integration has become increasingly urgent [2]. Still, effective policymaking in this area is obstructed by the lack of rudimentary knowledge about how immigrant integration responds to an increase in immigration.

To this end, in a recent work [3] a new approach for studying key-integration quantifiers, based on methods, models, and ideas from statistical physics, was proposed. The theory describes and predicts how typical integration quantifiers change when the density of migrants increases. The results predicted a linear growth for the averages of economic quantifiers like permanent and temporary jobs given to immigrant, and a square root growth for the averages of social quantifiers like mixed marriages and newborns to mixed couples. This framework is a powerful tool for the policy makers that are interested in assessing and evaluating integration progresses at the *municipality level*.

To deal with the phenomena at *municipality level* we use here a different theoretical framework based on

the theory and techniques of continuous random walks [13, 18]. The approach developed in [3], based on a full micro-macro statistical mechanics theory, revealed in fact a high efficacy to forecast average values. However, since the developed model doesn't have yet an exact solution, its related phase space picture is not fully disclosed and doesn't cover yet the structure of the fluctuations around the mean values. The random walk approach that we follow here instead, based on a meso-macro stochastic process, has the advantage to allow for a full analytical control of both mean values and fluctuations.

We consider classical quantifiers of integration such as the fraction of all temporary and permanent labor contracts given to immigrants, the fraction of marriages with spouses of mixed origin (native and immigrant), and the fraction of newborns with parents of mixed origin. The evolution of these quantifiers versus the percentage of migrants inside the host country is "locally erratic", that is, when looked at a fine level of resolution such as the municipality, it can be thought of as a *random walk* where the time change is represented by the change of migrant density in the municipality, and the integration quantifier – playing the role of the space variable – changes according to suitable probability distributions defining the stochastic process. Instead of obtaining the evolution of averages via statistical mechanics, with this approach the evolution of averages are here the result of averaging over the whole ensemble of municipalities, i.e., averaging over all the random walks.

From a sociological perspective, the evolution of the quantifiers, with respect to the density of immigrants, is, in fact, a random process whose stochasticity may depend on several exogenous factors driving immigration: fluctuations in the ratio between work demand and work request in the host country [2]), or "biases" resulting from (for example) *push-pull* factors [2] or different types of network induced migration outcomes [4–6]. However our aim here is not to explain or disentangle these mechanisms, but rather to look at the evolution of quantifiers as a combined effect of a "drift" in the presence of some "noise" regardless of its source/origin. To this task we use random walk theory: the latter constitutes the prototype of stochastic process, and, at the same time, the basic model of diffusion phenomena and non-deterministic motion. Indeed, applications can be found in the study of, for example, transport in disordered media (e.g., [7]), anomalous relaxation in polymer chains (see e.g., [8]), financial markets (see e.g., [9]), quantitative analysis in sports (see e.g., [10]).

Using stochastic process theory allows to get a mesoscopic description of the integration quantifiers behavior and to addresses questions such as whether these socio-economic metrics are determined by memory-less stochastic processes or by processes with long-time correlations. Moreover, this framework allow us to analyze rare events and non-Markovian quantities which are important determinants for planning, in so far they are key tools for quantifying fluctuations. That is, we aim to provide efficient tools to help assessing the progress (or deficit) in integration as well as to generate strong predictions for extreme case scenarios at lower administrative levels such as *municipalities*, and thereby, through an interplay between statistical mechanics and stochastic processes, we broaden the scope of practical applications of the quantitative theory of immigrant integration as a whole. Typical questions begging an answer are for example: What is the worst/best case scenario in the two integration branches – social and economic integration – in a particular municipality if immigrant density changes from say 5 to 7 percent? And how does the effect magnitude of this change compare to the effect magnitude of an equivalent change at the national level, i.e., average change, or in a similar/dissimilar municipality? In other words, through first-passage-time and maximum-span techniques, we obtain estimates for the expected value of immigrant density for which a particular integration quantifier – say, the share of immigrant workers or the number of mixed marriages – reaches a given threshold above which new policies, structures, services, facilities etc., have to be made available.

The work is organized as follows: first we describe the database and the procedures for data extraction (Sec. II), then we explain in details the mapping between the evolution of social quantifier and of a random walk (Sec. III and IV) and we report the related results (Sec. V). Finally, we discuss how such outcomes may be exploited to more effectively set up multiethnic plans and immi-

gration policies in general (Sec. VI).

II. DATA DESCRIPTION, ANALYSIS AND ELABORATION

Data considered here refer to quarterly observations during the period 1999 to 2010. It is drawn from Spain's Continuous Sample of Employment Histories (the so called *Muestra Continua de Vidas Laborales* or MCVL) [19] and from the local offices of Vital Records and Statistics across Spain (Registro Civil) [20]. The former provides detailed data on labor contracts, and the latter provides detailed data on spouses and parents to newborns. Information on the municipalities immigration density are drawn from the Municipal population registers [21]. A unique feature of the Spanish data is that all three data sources include also so called undocumented immigrants, that is, immigrants that lack a residence permit. Undocumented immigrants are usually not included in official statistical sources. However, their assimilation within the immigrant population is often significant and excluding them would underestimate the true size of the immigrant population as well as the frequency of the socio-economic events used to measure integration.

Because "municipality" is the lowest administrative level for which data on density is available, the individual data on mixed events is aggregated to the level of municipality. From these datasets, for each municipality [22] we obtain quarterly time series for the following quantities:

$$J_p = \frac{\text{\#permanent contracts to immigrant}}{\text{\#permanent contracts}}, \quad (1)$$

$$J_t = \frac{\text{\#temporary contracts to immigrant}}{\text{\#temporary contracts}}, \quad (2)$$

$$M_m = \frac{\text{\#mixed marriages}}{\text{\#marriages}}, \quad (3)$$

$$B_m = \frac{\text{\#newborns with mixed parents}}{\text{\#newborns}}. \quad (4)$$

As explained below, by studying how the quantities in Eqs. 1-4 vary with the overall fraction of immigrants, we can unveil the growth law determining their evolution and based on this information make previsions.

In order to assess the evolution of the Immigrants-Natives system, a convenient quantity to use as control parameter is

$$\Gamma = N_{imm}N_{nat}/N^2 = \gamma(1 - \gamma), \quad (5)$$

where $\gamma = N_{imm}/N$ is the fraction of immigrants. Indeed, Γ provides an intensive measure of the cross-links existing among the communities of natives and of immigrants (however, for small values of γ , $\Gamma \sim \gamma$, hence we can roughly map the percentage of migrants with the time in our bridge). Moreover, differently from other

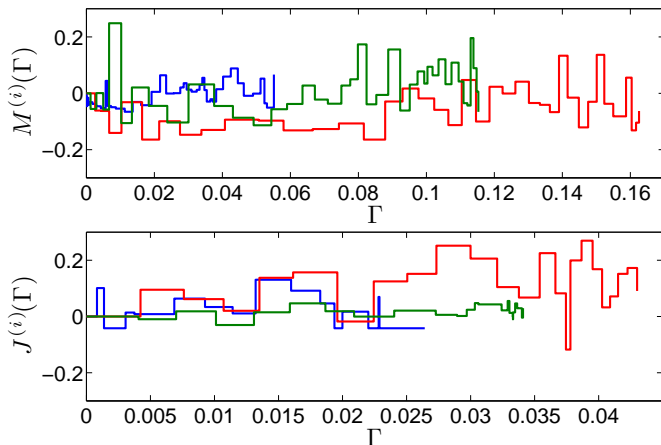


FIG. 1: Examples of paths for the quantifiers M_m (upper panel) and J_p (lower panel) shown as a function of Γ . Three different municipalities are depicted in different colors. These paths can be compared with a theoretical one depicted in Fig. 3 and related to a CTRW.

possible choices such as time, using Γ avoids any inaccuracy due to seasonality and allows to directly compare municipalities of different sizes (see also [3]).

Complete time series for data on labor contracts involve $\mathcal{M}_J = 124$ municipalities and consist of 2976 data entries over the period 2005–2010 which is sampled quarterly (i.e. overall 24 trimesters). Complete series for data on marriages and newborns involve $\mathcal{M}_M = 581$ municipalities and consist of 23240 data entries spanning the period 1999–2008 which is sampled quarterly (i.e. overall 40 trimesters).

Thus, for any municipality i , we consider five time series: one for Γ and one for each observable in Eqs. 1-4, hereafter denoted generically as $X^{(i)}$.

As Γ varies, each series $X^{(i)}$ determines a “path” in the related space and this point process can be looked at as a continuous-time random walk (CTRW) [23], where the time variable is given by Γ , while the space variable is given by $X^{(i)}$, see Fig. 1. This mapping is fully described in the next section.

Finally, in Fig. 2 we show the time series for $X^{(i)}$ and $\Gamma^{(i)}$ vs time (in units of trimesters) to highlight the different shape of paths.

A. Telegraphic introduction on CTRWs

A CTRW process can be depicted as a dynamical point (to fix ideas embedded in a one-dimensional space, as here we need such a case only), which occupies a position $r(t)$ at time t (see also Fig. 3). Let us suppose that the point starts on the origin, that is $r(0) = 0$. Then, it stays fixed to its position until time t_1 , when it jumps to ξ_1 , where it waits until time $t_2 > t_1$, when it jumps to a new location $\xi_1 + \xi_2$, and so on. The series $\{t_1, t_2, \dots\}$ defines

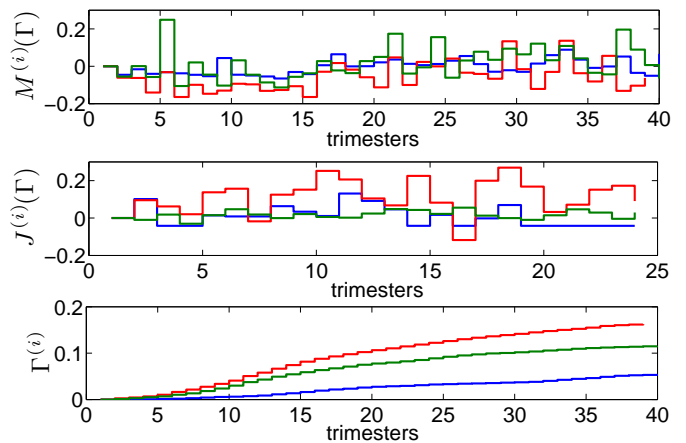


FIG. 2: Examples of paths for the quantifiers M_m (upper panel) and J_p (lower panel) shown as a function of time (1 unit = 1 trimester). Three different municipalities (the same as in Fig. 1) are depicted in different colors. Notice that for marriages seasonality effects emerge: during summer months marriages are more frequent.

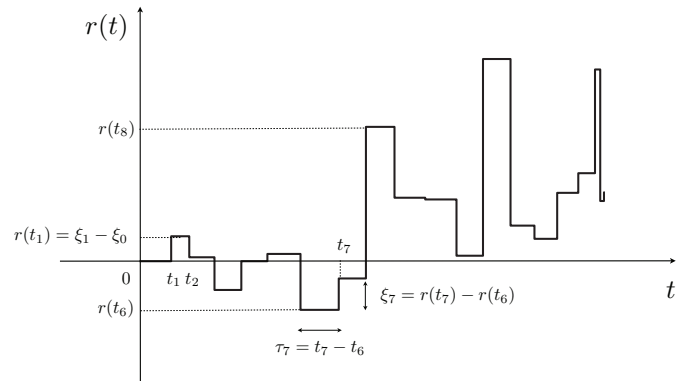


FIG. 3: Example of path realized by a CTRW whose step widths and waiting times are extracted from the distributions given by Eqs. 28 and 30, respectively, and with parameters consistent with those found experimentally (see Tab. II).

the times of jumping events. The times $\tau_1 = t_1 - 0$, $\tau_2 = t_2 - t_1$ etc. are called waiting times.

The waiting times $\{\tau_i\}$ and the width of the instantaneous jumps $\{\xi_i\}$ are continuous random variables extracted from the distribution $\psi(\xi, \tau)$. The latter determines the long-time properties of the walk: a diverging average waiting time typically corresponds to sub-diffusive behaviors, while a diverging variance for jump widths typically corresponds to super-diffusive behaviors.

In particular, for the so-called decoupled continuous random walk (namely where the distribution $\psi(\xi, \tau)$ factorizes into $\psi(\xi, \tau) = f(\xi)\psi(\tau)$), the waiting times and the instantaneous displacements are mutually independent (identically) distributed random variables.

The position r of the particle at the k -th jump, that is at time t_k , is given by the sum $r(t_k) = \sum_{i=1}^k \xi_i$. Getting

$r(t)$, namely a direct dependence on t , requires the introduction of the random variable $n(t)$, representing the number of steps m performed up to time t and defined by $n(t) = \max\{m : t_m \leq t\}$, in such a way that

$$r(t) = \sum_{i=1}^{n(t)} \xi_i. \quad (6)$$

The expected value $\overline{r}(t)$ of the displacement can be derived from the probability distributions for the waiting time and for the step length. In fact, focusing on the decoupled case [24], we can define $\overline{\xi} = \int \xi f(\xi) d\xi$ and $\overline{\tau} = \int \tau \psi(\tau) d\tau$, whereby, as long as $\overline{\tau}$ is finite, one can show that, in the limit of large t [11]

$$\overline{r}(t) \sim \overline{\xi} \frac{t}{\overline{\tau}}. \quad (7)$$

Thus, if there is no net drift ($\overline{\xi} = 0$), the average displacement is zero and one usually looks at the mean square displacement which turns out to scale as $\overline{r^2}(t) \sim \overline{\xi^2} t / \overline{\tau}$, and the purely diffusive limit can be recovered.

On the other hand, in the presence of a net drift ($\overline{\xi} \neq 0$), the mean displacement can also be expressed in terms of the mean number of steps $\overline{n}(t)$ performed up to time t as (see e.g., [11, 12])

$$\overline{r}(t) = \overline{n}(t) \cdot \overline{\xi}, \quad (8)$$

and, accordingly, $\overline{r^2}(t) \sim \overline{r}(t)^2$ [11, 12]. From Eq. 8, one can see that if the average time diverges or displays any anomalous behavior, the biased motion turns out to be anomalous as well.

Of course, the definitions given here can be extended to a geometrical space with arbitrary topology [13].

Despite this random walk process is, by definition, Markovian, one can also introduce non-Markovian related quantities such as the mean-first passage time \tilde{t} and the maximum span \tilde{r} , [14].

The mean-first passage time represents the mean time taken by a random walk to first reach a (fixed) point placed at a given initial distance r . Its dependence on r qualitatively depends on the kind of diffusion realized, in particular:

$$\tilde{t} \sim r^2, \text{ for pure diffusion} \quad (9)$$

$$\tilde{t} \sim r, \text{ for biased diffusion.} \quad (10)$$

The maximum span represents the farthest distance ever reached by a random walk up to time t . Again, the functional form of \tilde{r} as a function of t depends on the kind of diffusion realized:

$$\tilde{r} \sim \sqrt{t} \text{ for pure diffusion} \quad (11)$$

$$\tilde{r} \sim t, \text{ for biased diffusion.} \quad (12)$$

These relatively simple laws stem from the peculiarity of the one-dimensional structure. In general, the behavior of \tilde{t} and \tilde{r} functionally depends on the underlying topology.

Indeed, due to their non-Markovian nature, estimating such quantities may be rather tricky, yet they are intensively studied as they provide useful information and play an important role in many real situations (e.g. transport in disordered media, neuron firing, spread of diseases and target search processes [13, 15, 16]).

To summarize, the CTRW is a stochastic model for which $\psi(\tau)$ and $f(\xi)$ serve as input functions. The output is provided by the temporal series $\{t_1, t_2, \dots\}$ and $\{r_1, r_2, \dots\}$ from which quantities such as mean squared displacement, mean first-passage time, etc. can be calculated.

In the next section, the jump widths ξ_i 's as well as the positions $r(t)$ will assume different meanings (i.e., number of mixed marriages, of newborns from mixed couples, of temporary/permanent contracts to immigrants) according to the specific quantifier addressed.

III. THE MAPPING IN A NUTSHELL

Let us denote with $X^{(i)}$ a generic quantifier (i.e., the number of mixed marriages, of newborns from mixed couples, of temporary/permanent contracts to immigrants), where i specifies the municipality. According to the quantifier considered i is bounded by \mathcal{M}_J or by \mathcal{M}_M .

Therefore, we have the time series

$$\{X_1^{(i)}, X_2^{(i)}, \dots, X_{\mathcal{T}}^{(i)}\}, \quad (13)$$

$$\{\Gamma_1^{(i)}, \Gamma_2^{(i)}, \dots, \Gamma_{\mathcal{T}}^{(i)}\}, \quad (14)$$

where $X_n^{(i)}$ and $\Gamma_n^{(i)}$ are the values of the quantifier and of the number of cross-links at the n -th trimester and \mathcal{T} is bounded by the overall number of trimesters over which measures have been taken (i.e., 24 for job quantifiers and 40 for family quantifiers).

For a (one-dimensional) CTRW of \mathcal{T} steps, defined by the two series

$$\{\xi_1, \xi_2, \dots, \xi_{\mathcal{T}}\}, \quad (15)$$

$$\{t_1, t_2, \dots, t_{\mathcal{T}}\}, \quad (16)$$

where ξ_n is the jump width and t_n is time when the n -th step occurs, we recall that the position $r(t)$ of a walker at time t is obtained by $r(t) = \sum_{j=1}^{n(t)} \xi_j$, where $n(t)$ is the number of steps performed up to time t .

Analogously, we can state that, for the i -th municipality, the value of the quantifier $X^{(i)}(\Gamma)$ corresponding to degree of cross-link Γ is

$$X^{(i)}(\Gamma) = \sum_{j=1}^{n^{(i)}(\Gamma)} \Delta X_j^{(i)}, \quad (17)$$

where $\Delta X_j^{(i)} = X_{j+1}^{(i)} - X_j^{(i)}$ and $n^{(i)}(\Gamma)$ is the latest trimester for which $\Gamma_j^{(i)} < \Gamma$.

Therefore, we can look at the set of \mathcal{M} municipalities as a set of \mathcal{M} random walks. Actually, before proceeding, a couple of remarks are in order.

In principle, Γ and X are bounded by 1, yet, the number of immigrants corresponds to a small fraction of the overall population in such a way that $\Gamma, X \ll 1$ and we can neglect boundaries [25].

Moreover, Γ and X are not continuous variables as there exists an intrinsic unit given by $1/\text{\#number of marriages}$, $1/\text{\#number of newborns}$ and $1/\text{\#number of contracts}$, representing our experimental sensitivity. However, such a unit is in general much smaller than the quantities measured which can therefore be considered as continuous.

Therefore, we can treat the set of \mathcal{M} municipalities as a set of \mathcal{M} random walks, for which we can build the following ensemble average:

$$\langle X(\Gamma) \rangle \equiv \frac{1}{\mathcal{M}} \sum_{i=1}^{\mathcal{M}} X^{(i)}(\Gamma). \quad (18)$$

Similarly, for the average square distance covered

$$\langle X^2(\Gamma) \rangle \equiv \frac{1}{\mathcal{M}} \sum_{i=1}^{\mathcal{M}} [X^{(i)}(\Gamma)]^2. \quad (19)$$

The progression of the quantifiers $\langle X(\Gamma) \rangle$ averaged over the whole set of municipalities, that is to say, the average displacement of the related CTRW, is shown in Fig. 4, where fits evidence the following behaviors

$$\langle J_t(\Gamma) \rangle \sim \Gamma, \quad (20)$$

$$\langle J_p(\Gamma) \rangle \sim \Gamma, \quad (21)$$

$$\langle M_m(\Gamma) \rangle \sim \sqrt{\Gamma}, \quad (22)$$

$$\langle B_m(\Gamma) \rangle \sim \sqrt{\Gamma}. \quad (23)$$

perfectly consistent with those outlined in [3], despite the procedure for their derivation is conceptually different; this confers robustness to the above results.

To summarize, in our random-walk picture for the time evolution of the social quantifier X , in each municipality the quantifier starts from zero and, for a given variation of the related immigrant percentage Γ , the quantifier increases or decreases until the path ends. The trajectory of X versus Γ qualitatively resembles the position of a CTRW as a function of time (see Figs. 1 and 3).

In the next section we analyze the CTRWs associated to the quantifiers and try to get a *microscopic* perspective for the origin of these laws. Such a perspective will allow to speculate about possible effects and to make crucial forecasts.

IV. FORMALIZING THE MAPPING

We first check that the CTRWs corresponding to J_p , J_t , M_m and B_m are decoupled, that is, the related probability distributions $\psi(\Delta X, \Delta \Gamma)$ for the generic increments

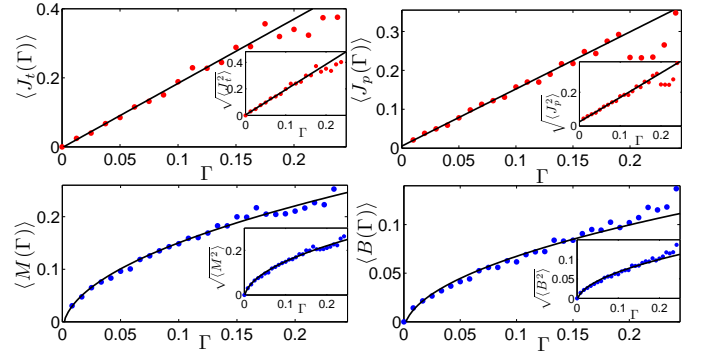


FIG. 4: “Mean displacement” (main figures) and “mean square displacement” (insets) versus “time” for the CTRWs associated to J_t (panel a), J_p (panel b), M_m (panel c) and B_m (panel d). Data available were binned over Γ and averaged over the set of \mathcal{M} municipalities; the resulting values (\bullet) and the related best fit (solid line) are shown. In particular, for family quantifiers we fitted by the law $r = p_1\sqrt{t} + p_2$, while for job quantifiers we used the law $r = p_3t + p_4$; best fit coefficients are summarized in Tab. I. In general, the goodness-of-fit R^2 ranges between 0.97 and 0.99. Notice that $\sqrt{\langle X^2(\Gamma) \rangle} \sim \langle X(\Gamma) \rangle$ suggests the presence of a drift [11].

Quantifier X	p_1	p_2
$\langle M_m \rangle$	0.54 ± 0.02	-0.019 ± 0.009
$\sqrt{\langle M_m^2 \rangle}$	0.57 ± 0.03	0.007 ± 0.06
$\langle B_m \rangle$	0.25 ± 0.01	-0.010 ± 0.009
$\sqrt{\langle B_m^2 \rangle}$	0.287 ± 0.002	-0.007 ± 0.004
Quantifier X	p_3	p_4
$\langle J_t \rangle$	1.9 ± 0.1	-0.003 ± 0.001
$\sqrt{\langle J_t^2 \rangle}$	1.9 ± 0.1	0.003 ± 0.001
$\langle J_p \rangle$	1.47 ± 0.06	0.005 ± 0.003
$\sqrt{\langle J_p^2 \rangle}$	1.45 ± 0.07	0.025 ± 0.008

TABLE I: Best-fit coefficient related to plots shown in Fig. 4.

ΔX and $\Delta \Gamma$ can be factorized into $f(\Delta X)\psi(\Delta \Gamma)$: this is achieved through direct inspection of the scatter plots reported in Fig. 5.

Thus, we can proceed by studying separately $f(\Delta X)$ and $\psi(\Delta \Gamma)$. We recall that such distributions provide qualitative information about the diffusive behaviors of the walks associated to our quantifiers, that is, on their time progress. Moreover, from $f(\Delta X)$ and $\psi(\Delta \Gamma)$, we are able to derive the expectation values

$$\overline{\Delta X} = \int \Delta X f(\Delta X) d\Delta X, \quad (24)$$

$$\overline{\Delta \Gamma} = \int \Delta \Gamma \psi(\Delta \Gamma) d\Delta \Gamma, \quad (25)$$

which play as the expected jump length and as the expected waiting time respectively. Analogously, we can derive $n(\Gamma)$ which plays as the expected number of steps performed up to “time” Γ , that is

$$\overline{n(\Gamma)} = \sum_n n Q(n|\Gamma), \quad (26)$$

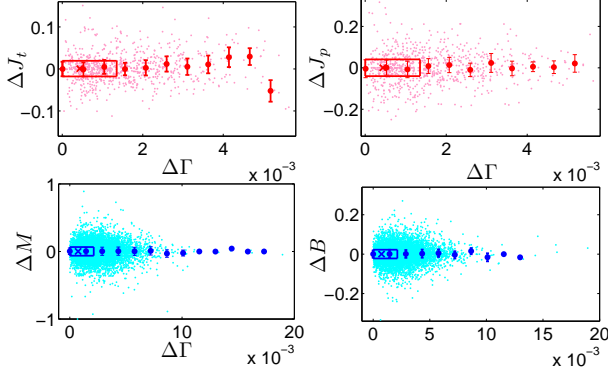


FIG. 5: These scatter plots evidence the existence of any correlation between the “waiting times” $\Delta\Gamma$ and the “jump width” ΔM (panel a), ΔJ (panel b), ΔJ_t (panel c), ΔB (panel d): each point represents the increments ΔX_n versus $\Delta\Gamma_n$; all \mathcal{T} steps and the whole set of municipalities are considered. The clouds of data are uniform and do not reveal any special trend. Binned spots evidence the possible values of increments $\Delta\Gamma$ and for each bin we calculated the average of the related increments ΔX_n ; the related standard deviation are also depicted. Notice that such averages are basically constant (at least within the error) with respect to $\Delta\Gamma$ and this allows to derive that no clear correlation emerges.

where $Q(n|\Gamma)$ is the probability that $\sum_j^n \Delta\Gamma_j$ is smaller than Γ , but $\sum_j^{n+1} \Delta\Gamma_j$ is larger than Γ .

From these quantities, one finally has (see e.g., [11, 12])

$$\overline{X(\Gamma)} = \overline{n(\Gamma)} \cdot \overline{\Delta X}. \quad (27)$$

Of course, the expectation $\overline{X(\Gamma)}$ and the ensemble average $\langle X(\Gamma) \rangle$ ought to be consistent (as checked in the next section). This ensures the ergodicity of the system and will allow us to exploit the analytical results derived starting from the probability distribution functions also for our “time” series.

A. Step width and Waiting time distributions

Let us start with the distribution for the “step lengths” $f(\Delta X)$. In Fig. 6 we show the histogram for the increments ΔJ_t , ΔM , ΔJ_p and ΔB obtained from experimental data. In all cases the symmetric, centered exponential distribution

$$f(\Delta X) = \lambda e^{-\lambda|\Delta X|}, \quad (28)$$

provides an excellent fit. An exponential distribution for step lengths ensures that the related CTRW does not exhibit any super-diffusive feature as the central limit theorem is fulfilled.

Now, the fit coefficient λ depends on the quantifier considered and it is directly related to the expected value by $\lambda_X^{-1} = \overline{\Delta X}$. Results are collected in Tab. II, where a comparison with the experimental average values $\langle |\Delta X| \rangle$ and $\langle \Delta X \rangle$ is also provided.

Quantifier X	λ_X^{-1}	$\langle \Delta X \rangle$	$\langle \Delta X \rangle$
J_t	0.031 ± 0.002	0.03	0.003
J_p	0.058 ± 0.003	0.06	0.003
M_m	0.079 ± 0.002	0.08	0.003
B_m	0.035 ± 0.001	0.03	0.001

TABLE II: The second column contains the best-fit coefficients obtained by fitting, according to Eq. 28, the probability distribution function of the displacements ΔX shown in Fig. 6, while the third and fourth columns contain the related average values, where the average is performed on raw data over all municipalities. Being the support of the exponential distribution positive, λ_X^{-1} has to be compared with $\langle |\Delta X| \rangle$. Moreover, we checked that the absolute error on $\langle |\Delta X| \rangle$ is approximately equal to $\langle |\Delta X| \rangle$ itself, as expected from an exponentially-distributed variable. Notice that the average displacement $\langle \Delta X \rangle$ in a single step is positive for any quantifier.

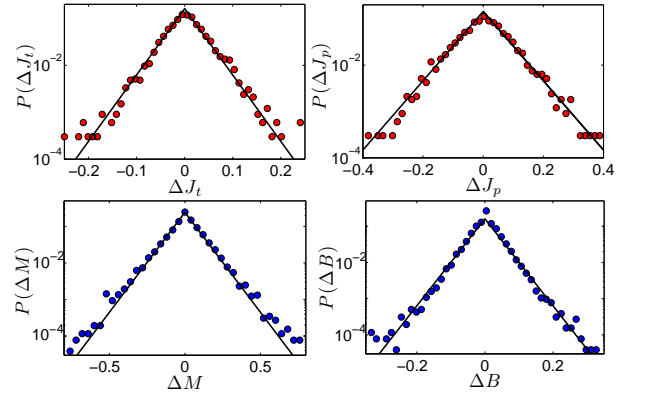


FIG. 6: Distributions $f(\Delta M)$ (panel a), $f(\Delta B)$ (panel b), $f(\Delta J_t)$ (panel c) and $f(\Delta J_p)$ (panel d) measured from experimental data, without distinguishing between municipalities, that is, we merged the increments pertaining to the whole ensemble of walks and we built a unique histogram. Notice the semi-logarithmic scale plot. Data (\bullet) are fitted by using Eq. 28 (solid line); best-fit coefficients and averages on raw data are collected in Tab. II

The goodness of the fit is corroborated by the fact that λ_X^{-1} and $\langle |\Delta X| \rangle$ coincide within the error. However, looking at $\langle \Delta X \rangle$ we report a slight deviation: while one would expect a null average value due to the centrality of the distribution, the average is systematically positive for all quantifiers and this implies that, as Γ increases, X is more likely to grow rather than to decrease. In the random-walk picture, this can be interpreted as the presence of a drift which biases the motion of the walker.

Let us now move to the distribution for the “waiting times” $\Delta\Gamma$.

In Fig. 7 we show the histogram for the increments $\Delta\Gamma$ obtained from experimental data related to the time period and to the municipalities considered. Interestingly, here qualitative differences emerge between the job quantifiers, i.e. J_t and J_p , and the family quantifiers, i.e. M_m

and B_m .

Before proceeding it is worth stressing that for job quantifiers and family quantifiers the time along which sampling has been performed is not exactly the same, being, respectively, 2005-2010 and 1999-2008 (of course, the consistency between the related time series has been checked for the overlapping period [3]). Now, in order to ensure that the qualitative differences reported do not stem from different time interval, but are intrinsic, we repeated the analysis shown in Fig. 7 by restricting only to the common time lapse 2005-2008 and, indeed, we checked the robustness of the result.

In fact, calling ψ_F and ψ_J the distributions for family and job quantifiers respectively, the reason for their intrinsic difference can be depicted in the way mapping between quantifier evolution and random-walks has been fixed. In particular, there exist trimesters i for which a growth in the number of immigrants is reported, i.e. $\Gamma_i - \Gamma_{i-1} > 0$, but no change in the quantifier X considered occurs, i.e. $\Delta X_i = 0$. In such cases the two trimesters behave as practically merged as the overall waiting time gets $\Gamma_{i+1} - \Gamma_{i-1}$. This concept can be repeated iteratively until each step of the walk actually corresponds to a true displacement. Thus, as one can see from Fig. 7, such merging are more frequent for family quantifiers in such a way that the related waiting times display a larger range. Otherwise stated, the integration of immigrants within the market is more direct: as long as new immigrants arrive, a fraction of them get a job, either permanent or temporary. Conversely, the integration of immigrants from a familiar perspective is more complex and does not follow a prescribed pattern: not surprisingly, the arrival of new immigrants does not necessarily correspond to integration when considering these quantifiers. This is consistent with the results in [3], where from a different perspective, it is shown that the qualitative difference between the laws $M_m(\Gamma)$, $B_m(\Gamma)$ and $J_t(\Gamma)$, $J_p(\Gamma)$ is due to a different degree of interaction among agents in the two different scenario (families and jobs).

It is worth stressing that such effect is not directly imputable to the seasonality of marriages; this can be seen, for instance, from the fact that for newborns the same effect emerges as well, but their time series do not display any seasonality.

Let us now analyze in more details the waiting time distributions.

For family quantifiers the distribution $\psi_F(\Delta\Gamma)$ fitting the experimental histogram is a log-normal distribution

$$\psi_F(\Delta\Gamma) = \frac{1}{\Delta\Gamma\sqrt{2\pi}\sigma} \exp - \frac{(\log \Delta\Gamma - \mu)^2}{2\sigma^2}, \quad (29)$$

for which the average value is expected to be $\overline{\Delta X} = e^{\mu+\sigma^2}$. As for jobs, the best fit is provided by a half-normal distribution

$$\psi_J(\Delta\Gamma) = \frac{\sqrt{2}}{\sqrt{\pi}\sigma} \exp - \frac{(\Delta\Gamma - \mu)^2}{2\sigma^2}, \quad (30)$$

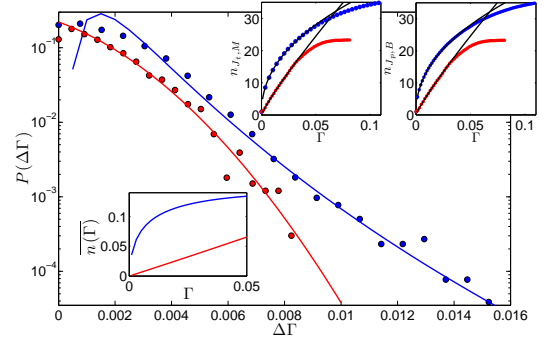


FIG. 7: Main plot: Histograms for $\Delta\Gamma$, derived from experimental data concerning marriages (blue symbols) and permanent jobs (red symbols), are shown and compared. Solid lines represent the best fit according to a lognormal distribution (see Eq. 29) and a half-Gaussian distribution (see Eq. 30), respectively. Fitting coefficients and related errors are reported in Tab. III. Notice that such histograms were derived without distinguishing between municipalities. Lower inset: average number of steps performed up to time Γ , calculated numerically from Eq. 29 (red line) and Eq. 30 (blue line), respectively. Upper insets: Average number $\langle n(\Gamma) \rangle$ of steps performed by the related random walker up to the fraction of immigrants Γ . Solid lines correspond to the law $y \sim x$ and $y \sim \sqrt{x}$, respectively and evidence qualitative different behaviors for marriages and jobs. This picture corroborates the validity of Eq. 8 with the ensemble average: $\langle r \rangle \sim \langle n \rangle \langle \Delta r \rangle$, which bridges the picture itself with Fig. 1. The fit is robust only up to relatively small values of Γ , then experimental averages are underestimated. This is due to the fact that the statistics is robust only for values of Γ which are reached by (almost) all walks. For larger values our averages are only an underestimate of the expected, effective mean value of n .

for which the average value is expected to be $\overline{\Delta X} = \mu$. Details on fitting coefficients and average values are all collected in Tab. III; notice that, in both cases, $\overline{\Delta X}$ turns out to be comparable with the ensemble average $\langle \Delta\Gamma \rangle$.

Thus, although both ψ_J and ψ_F fulfill the central limit theorem and display a finite mean, the latter displays a long tail so that we expect that the growth for family quantifiers may be slowed down.

In particular, we expect such slowing down to be more evident at “short times”, namely for small values of Γ . This can be seen intuitively: for family quantifiers waiting times are more broadly distributed in such a way that for relatively small values of Γ it is likely that the number n of steps performed is rather small, that is, smaller than the mean-field expectation value $\Gamma/\langle \Delta\Gamma \rangle$.

Now, given ψ_J and ψ_F , we can derive the number of steps performed up to time Γ , exploiting the properties of Laplace transforms (see e.g., [11, 12]). Examples of numerical results of these calculations are shown in the lower inset of Fig. 7: the difference between the two cases is striking.

In order to check this point we measure directly on raw data the average number $\langle n(\Gamma) \rangle$ of steps performed before reaching the time Γ (see Fig. 7). Indeed, for jobs

Γ	μ	σ^2	$\overline{\Delta\Gamma}$	$\langle\Delta\Gamma\rangle$
Job	$(1.2 \pm 0.2) \cdot 10^{-3}$	$(6.7 \pm 0.6) \cdot 10^{-6}$	$(2.0 \pm 0.2) \cdot 10^{-3}$	$(1.7 \pm 0.2) \cdot 10^{-3}$
Family	-6.6 ± 0.9	0.32 ± 0.04	$(1.7 \pm 0.3) \cdot 10^{-3}$	$(1.9 \pm 0.3) \cdot 10^{-3}$

TABLE III: Best-fit coefficients obtained by fitting the probability distribution function of the “waiting time” $\Delta\Gamma$ shown in Fig. 7 according to Eqs. 29 and 30. The relative error on fit coefficients ranges between 10% and 20%. Within the error there is perfect consistency between the average values $\overline{\Delta X}$ and $\langle\Delta X\rangle$, as well as between the variance of such distributions and the variance on the related raw data. Here we report only data for marriages and permanent jobs; for newborns and temporary jobs analogous analysis evidence only slight quantitative changes.

we find a roughly linear growth, i.e. $\langle n(\Gamma) \rangle \sim \Gamma$, while for marriages and births we find a slower growth, i.e. $\langle n(\Gamma) \rangle \sim \sqrt{\Gamma}$.

Such a qualitative difference, together with Eq. 8, immediately explains the results of Eqs. 20-23.

Summarizing, both processes display a non-null positive drift, i.e. $\Delta X > 0$, yet the resulting behaviors are qualitatively different over the time window considered. Such a difference ultimately stems from deep differences in the waiting times: a broader distribution for $\Delta\Gamma$ occurs in the case of family quantifiers and the related random walks may experience rather long waiting times, although the jump widths remain narrowly distributed. The net result is just a slowing down in the progress of the quantifier.

Conversely, as for job, both ΔX and $\Delta\Gamma$ are narrowly distributed so that at each trimester we do not expect strong variations in the fraction of new immigrants getting a job.

Such a difference suggests an intuitive motivation, namely that the mechanisms underlying the emergence of mixed marriages are more complex and may be subjected to mutual interaction among individuals. This is perfectly consistent with the statistical-mechanics description of the phenomenon provided in [3].

V. FIRST PREDICTIVE OUTCOMES FOR SOCIAL PLANNERS

We now turn to the theory’s predictive capacity. The aim is to present concrete instruments directed to aid policy makers at the municipal level in their work to accommodate and plan for further immigration. We focus on two well-known observables: the (mean) first passage time, and the (mean) maximum walk span.

A. Mean first-passage time

Mean first-passage-time quantities have been extensively investigated in a number of different fields, ranging from chemical kinetics to finance, as they provide an estimate for the average time at which a given stochastic event is triggered [15, 16].

Given the process $X(\Gamma)$ we calculate the value $\tilde{\Gamma}(x)$ at

which the quantifier reach a certain threshold x . In order to evaluate the typical value of $\tilde{\Gamma}(x)$ we perform an average over the ensemble of walks, that is

$$\langle \tilde{\Gamma}(x) \rangle = \frac{1}{\mathcal{M}} \sum_{i=1}^{\mathcal{M}} \tilde{\Gamma}^{(i)}(x). \quad (31)$$

The quantity $\langle \tilde{\Gamma}(x) \rangle$ allows predictions about the consequences additional immigration have on integration and when a integration threshold is likely to be reached. For instance, let us say that when a integration quantifier reach the threshold x , some integration policies, activities, or services must be activated (e.g. concerning public education, public health, etc). Then, as Γ approaches $\langle \tilde{\Gamma}(x) \rangle$ local projects and plans need to be activated.

In Fig. 8 we show the mean-first passage time for the quantifiers considered in this work as a function of X . The mean first-passage time is especially useful for policies plans and service that are coupled with a concrete “discrete” integration target, and when we need to know the expected time when the politically defined threshold is reached, and activation of the plans are being called for.

For example, we could ask at which value of Γ (which is the percentage of migrants) we expect that the amount of newborns from mixed parents reaches the threshold of 10%. By simply looking at the behavior of $\langle X(\Gamma) \rangle$, by inverting, we would get $\Gamma \sim 0.2$. However, due to huge fluctuations (hence in some peculiar municipalities), the threshold of 10% can be reached much earlier, as the first passage time, returns a value $\Gamma \sim 0.04$. Hence planning based on average evolutions only may underestimate reality by a factor rendering planning and resource allocation extremely ineffective.

B. Walk span

The walk span represents the largest point reached by the walker up to a given time. That is, the largest value \tilde{X} reached by X up to Γ . More precisely, we say that for the i -th walk, at the k -th step, the span is $\tilde{X}^{(i)}$ if $X(n)^{(i)} < \tilde{X}^{(i)}(k), \forall n \leq k$. Again, in order to evaluate the typical value of $\tilde{X}(k)$ we perform an average over the ensemble of walks, that is

$$\langle \tilde{X}(k) \rangle = \frac{1}{\mathcal{M}} \sum_{i=1}^{\mathcal{M}} \tilde{X}^{(i)}(k). \quad (32)$$

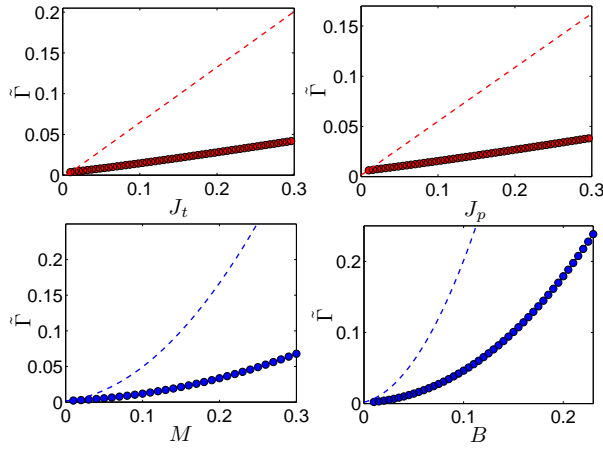


FIG. 8: Mean time $\tilde{\Gamma}$ to first reach a given value of J_t (panel a), of J_p (panel b), of M (panel c), of B (panel d). Experimental data (\bullet) are obtained by first getting the mean number of steps to first reach the distance X and then by inverting through $n(\Gamma)$ (see Fig. 7). Solid lines are best fits given by $y = p'_1\Gamma + p'_2$ (upper panels) and by $y = p'_3\Gamma^2 + p'_4$ (lower panels), being $p'_1 = 0.13 \pm 0.02$, $p'_2 = 0.0014 \pm 0.0005$ for J_p , $p'_1 = 0.11 \pm 0.02$, $p'_2 = 0.0045 \pm 0.0001$ for J_t and $p'_3 = 0.70 \pm 0.01$, $p'_4 = 0.0047 \pm 0.0003$ for M_m , $p'_3 = 4.54 \pm 0.03$, $p'_4 = 0.00044 \pm 0.0002$ for B_m . These results are compared with the related $\Gamma(X)$ (dashed line) derived from results shown in Fig. 4; see also data in Tab. I for comparison.

The average walk span provides information the capacity to integrate further immigration. In fact, in organizing local integration policies and make appropriate priority decisions among different integration initiatives, one is interested in the span of, say, the number of children, or the number of immigrants with permanent jobs, rather than in their average number as the latter may lead to dramatic over- and underestimations.

In Fig. 9 we show the span of the quantifiers considered in this work as a function of Γ . We notice that the qualitative differences already evidenced for $\langle X(\Gamma) \rangle$ are robust and, the span for marriages and births grows like $\sqrt{\Gamma}$, while the span for temporary and permanent jobs grows like Γ . The persistence of such behaviors is consistent with the fact that such random walks display distributions for waiting time and step width having finite average and variance. For instance, for a simple random walk on a line the span grows in time like \sqrt{t} , while in the presence of a drift one has a linear law t [13].

VI. CONCLUSIONS

Theoretical models, originally developed to solve physical problems, are increasingly being used to study social phenomena. Statistical mechanics and stochastic process theory are particularly well suited for this task, and have generated a novel quantitative understanding of the underlying complexity of social interactions. In this pa-

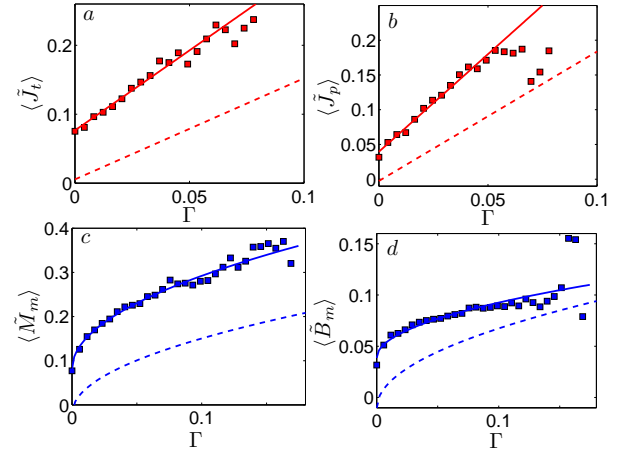


FIG. 9: Span of the walk for permanent jobs (panel a), for temporary jobs (panel b), for marriages (panel c) and for newborns (panel d) versus Γ . Solid lines are best fits given by $y = p'_1\Gamma + p'_2$ (upper panels) and by $y = p'_3\sqrt{\Gamma} + p'_4$ (lower panels), being $p'_1 = 2.3 \pm 0.2$, $p'_2 = 0.08 \pm 0.01$ for J_p , $p'_1 = 2.8 \pm 0.2$, $p'_2 = 0.04 \pm 0.01$ for J_t and $p'_3 = 0.8 \pm 0.1$, $p'_4 = 0.04 \pm 0.01$ for M_m , $p'_3 = 0.17 \pm 0.02$, $p'_4 = 0.04 \pm 0.01$ for B_m . These results are also compared with the curves $X(\Gamma)$ from Fig. 4 (dashed line); see also data in Tab. I for comparison.

per we focused on stochastic processes. We identified the random behavior of the four integration quantifiers with random walkers: each municipality draws a random walk in the quantifier-migrant's density plane. Averaging over all the municipalities then allowed to investigate the evolution of the quantifier averages, which are found to scale with the square root of the percentage of migrants for familiar quantifiers and linearly with the percentage of migrants for job quantifiers, in complete agreement with previous findings obtained through the statistical mechanical route [2]. We inferred the distributions of jumps and waiting times (which are found to be decoupled): while jump distributions are exponentially distributed for all the quantifiers, waiting time distributions depends on the context: social quantifiers have log-normal distributions for those times, while economic quantifiers display Gaussian distributions.

This difference has a simple explanation. While there is a correlation, even on a short timescale between the last-arrived migrant and that migrants incorporation into the labor market (in order to sustain), the same is not true for marriages or newborns. Clearly correlation is likely to be negligible between the last arrived migrant and a mixed marriage or birth event (i.e., it is unlikely that the arriving migrant and the one, say, marrying a native are the same person). This results in a stronger noise affecting our social quantifiers, which destroys the net drift and simple diffusion is the only survivor. On the contrary, driven by the necessity to work to survive, our economic quantifier display ballistic motion. Another motivation that contributes to the macroscopic differences resides in the much broader distribution of

jumps for the working quantifiers: The fat tail encoding for the long jumps in the working quantifiers implies a larger value of drift, that, coupled with much less noise – for the reasons just mentioned – lead to ballistic motion.

From a practical purpose, no power-law distributions are found. Hence, the Central Limit Theorem holds implying that the theory is suitable for generating predictions. To this end, we introduced two predictive non-Markovian tools: the “mean first passage time” and the “maximum span walk”. Using these tools it become possible to tackle questions that traditionally been answered using guesstimates in a more scientific way. For example, our predictive framework can easily produce forecast of the share of newborns with mixed parents following an increase in the share of immigrants from, say, 3 to 5 percent? We make two types of forecasts: first, we assess the evolution of the mean of this quantifier. The evolution is obtained evaluating from Figure 4 the average increment, which is roughly from $B(\Gamma) = 0.04$ to $B(\Gamma) = 0.05$. Second we assess the mean worst case by dealing with fluctuations. These fluctuations are obtained by extrapolating data from Figure 8, which gives a $\tilde{B}(\Gamma) \sim 0.08$, i.e. more than fifty percent higher than its average value. Although the investigated quantities are non-Markovian ($\langle \tilde{X}(\Gamma) \rangle$ and $\langle \tilde{T}(X) \rangle$) their behavior is still treatable: each of them can indeed be studied separately as a one-dimensional random walk also concerning the first passage time and the maximum span walk.

On a broader level, this work provides a concrete rig-

orous method for quantitative studies of social-science problems. The choice of immigrant integration is motivated by its prominent place in both the UE and the US political agendas. By uncovering the local variation pattern in the quantifiers we produced a scientific tool for anticipating the consequences of further immigration on local integration process. Information of this type has not been available in the past and constitutes great value for the development of immigration policies and multi-ethnic planning at the local level. However, while this work advances our knowledge on integration phenomena, other effects, like segregation phenomena, that may spontaneously develop in the host country has yet to be considered and incorporated into the theoretical framework developed here.

Acknowledgments

This work is supported by the FIRB grants RBFR08EKEV and RBFR10N90W.

EA and AB acknowledge also partial financial support by GNFM-(INdAM).

RS is grateful to the project Competition, Adaptation and Labour-Market Attainment of International Migrants in Europe (CALMA) granted by the VI National Plan for Scientific Research, Spanish Ministry of Economy and Competitiveness (CSO2012-38521), for partial financial support.

-
- [1] *European Commission: Handbook on Integration for policy-makers and practitioners* (Publications Office of the European Union, Luxembourg, 2010).
 - [2] S. Castles and M. J. Miller, *The Age of Migration - International Population Movements in the Modern World* (Pallgrave MacMillan, New York, 2009).
 - [3] A. Barra, P. Contucci, R. Sandell, and C. Vernia, *Sci. Rep. Nature* **4**, 4174 (2014).
 - [4] D. Massey and R. Zenteno, *Proc. Natl. Acad. Sci.* **96**, 5328 (1999).
 - [5] K. L. Wilson and A. Portes, *American Journal of Sociology* **86**, 295 (1980).
 - [6] R. Sandell, *Int. Migrat. Rev.* **49**, 971 (2012).
 - [7] E. Montroll and M. Schlesinger, *Nonequilibrium Phenomena II: From Stochastics to Hydrodynamics* (North-Holland, 1984).
 - [8] B. Hughes, E. Montroll, and M. Schlesinger, *J. Stat. Phys.* **28**, 111 (1982).
 - [9] J. Bouchaud and M. Potters, *Theory of Financial Risk and Derivative Pricing* (Cambridge University Press, 2003).
 - [10] A. Gabel and S. Redner, *Journal of Quantitative Analysis in Sports* **1416** (2012).
 - [11] R. Klages, G. Radons, and I. M. Sokolov, eds., *Anomalous Transport: Foundations and Applications* (Wiley-VCH, 2007).
 - [12] E. B. R. Metzler and J. Klafter, *Physica A* **266**, 343 (1999).
 - [13] G. Weiss, *Aspects and Applications of the Random Walk* (North-Holland, 1994).
 - [14] S. N. Majumdar, *Physica A* **389**, 4299 (2010).
 - [15] S. Redner, *A Guide to First-Passage Processes* (Cambridge, 2001).
 - [16] S. R. R. Metzler, G. Oshanin, ed., *First Passage Phenomena and Their Applications* (World Scientific, 2013).
 - [17] E. Montroll and G. Weiss, *J. Math. Phys.* **6**, 167 (1965).
 - [18] J. Klafter and I. Sokolov, *First Steps in Random Walks* (Oxford Press, 2011).
 - [19] It is an administrative data set with longitudinal information for a 4% non-stratified random sample of the population who are affiliated with Spain’s Social Security. We use data from the waves 2005 to 2010. The residence municipality is only disclosed if the population is larger than 40000.
 - [20] These data are compounded by the “National Statistical Agency” (INE). The residence municipality is only disclosed if the population is larger than 10000.
 - [21] More precisely, we use the size of the immigrant population and the native population in each municipality as reported in the 2001 Census as our baseline. Thereafter, based on the information contained in the “Statistics over residential variation in Spanish municipalities” and statistics on vital events (births and deaths) as elaborated by Spain’s “National Statistical Agency” (INE),

we estimate local immigrant densities for different points in time between 1999 and 2010.

- [22] Due to data protection, data on mixed marriages and newborns with mixed parents is only available for municipalities with a population larger than 10000. In addition, and due to data protection, municipality coding for the labor contract data is only available if the municipality's population exceeds 40000. However, about 85% of Spain's immigrants reside in the included municipalities.
- [23] The continuous time random walk (CTRW) was introduced by Montroll and Weiss [17]; see also [13, 18] for recent reviews and SI for a deeper description.
- [24] As we will show, this is the case recovered by our experimental data
- [25] Conversely, if boundaries can not be neglected the mapping could still be feasible but we should refer to the theory of random walks on finite chains

